

VISUAL QUESTION ANSWERING

Dr. Rhea Srinivas Department of Computer Science Jain University Bangalore, India rhea.sriniwas@jainuniversity.ac.in

S Rohith Kumar Computer Science and Engineering Jain University Bangalore, India seemakurthirohith@gmail.com R Navneeth Naidu Computer Science and Engineering Jain University Bangalore, India <u>navneeth.naidu03@gmail.com</u>

Santosh Nallala Computer Science and Engineering Jain University Bangalore, India santoshnallalarbp@gmail.com P Anil Kumar Computer Science and Engineering Jain University Bangalore, India anilnani8247@gmail.com

Mithun P Computer Science and Engineering Jain University Bangalore, India 21btrcs035@jainuniversity.ac.in

(VLP)Abstract Vision-Language Pre-Training significantly improves performance for a variety of multimodal tasks. However, existing models are often specialized in understanding or generation, which limits their versatility. Furthermore, trust in text data for large, loud web text remains the optimal approach for monitoring. To address these challenges, we propose VLX, a uniform VLP framework that distinguishes both vision languages and generation tasks. VLX introduces a new type of data optimization strategy. This strategy allows the generator to create high-quality synthetic training data, highlight the identifier noise, and allow the web to use the data records collected to more efficiently use the data records. Our framework achieves cutting-edge results with important benchmarks, including image text call (+3.1% average recall @1), visual answer questions (+2.0% accuracy), and multimodal capacitage (+2.5% cider). Additionally, VLX demonstrates the robust transferability of zero-shot transmissions to video language tasks without any additional tweaks. Publish codes, models and data records to promote future research.

I Introduction

Recent advancements in vision languages (vlp) have enhanced performance in multimodal tasks, but there are significant drawbacks: model rigidity two and unpredictable data. Endely models like clip and albet excel understand tasks like invocations, but they do run badly in text generation, while encoder decoder models like simvlm are better in generating, but are less effective at calling and limiting adaptability in different applications. Furthermore, state-of-the-art vlps are frequently based on image-text pairs found in web sources. In this pair, it hinders or significantly disrupts the learning process, and the data size does not completely eliminate this issue. In order to tackle these difficulties, we suggest blip (pretraining bootstrap language images), a uniform vlp framework that enhances both the model's structure and the quality of the data.

Architectural multimodal mixed (med) encoder decoder (med) works flexibly as an image encoder for calls and as an image-oriented decoder for generation, as well as a unimodal encoder. Equipped with contrasting learning, matching, and language modeling, med enables smooth transition between tasks. By incorporating this, our captions and filtering (capfilt) techniques enhance web data by generating synthetic sympathy, creating inferior infringement, and establishing a self-daily loop that improves data efficiency without manual labeling. Blip achieves groundbreaking results, concludes, questions and visual conversations. This demonstrates a substantial enhancement, with a +2.7% @1 improvement on call, and also showcases a remarkable generalization of zero-shot video language tasks. Ablation studies validate the significance of various synthetic caps and efficient filtering mechanisms. By integrating knowledge of vision languages with techniques for generating and reducing noise in data, blip establishes a solid base for a comprehensive multimodal system that can be openly shared with code, models, and data records to facilitate future research.

II Related Works

2.1 Vision-Language Models The field of vlp has experienced significant growth and advancement in recent years. Models such as clip, align, and albef employ contrastive learning objectives to align images and text in a shared embedding space. These models are exceptionally efficient in tasks such as zero-shot classification and imagetext retrieval, but they do not possess the ability to generate new content. Encoder-decoder models such as simvlm, blip, flamingo, and git have the capability to generate captions for images and answer questions in natural language through conditional text generation. Although versatile. these models often sacrifice retrieval performance and face increased computational complexity. Recent advancements also encompass unified frameworks that strive to strike a balance between comprehending and generating capabilities, but they often struggle to be



effectively implemented on a large scale due to training inefficiencies and limited adaptability to different data sets. Vlx fills this gap by providing a flexible hybrid architecture and utilizing automated data enhancement, pushing the limits of what unified vlp systems can accomplish.

2.2 Data Quality and Optimization The effectiveness of vlp models heavily relies on the quality of the training dataset. Datasets such as laion-400m, conceptual captions, and yfcc100m are commonly utilized because of their large size but are inherently prone to noise. In the past, attempts to enhance data quality have involved heuristic filtering, manual annotation, and the utilization of pre-trained language models to either re-rank or correct captions. Blip introduced the idea of improving caption quality by using synthetic generation. Unfortunately, these approaches either demand extra manual work or do not fully utilize the iterative nature of self-improvement. Vlx improves data optimization by implementing an autonomous feedback mechanism, allowing the model to refine its training dataset without requiring manual intervention. This repetitive cycle guarantees that the quality of the training data consistently enhances, establishing a strong basis for model convergence and generalization across various tasks and domains.

III Methodology

3.1 Unified Model Architecture VLX introduces a hybrid multimodal architecture that integrates the strengths of encoder-only and encoder-decoder models. The model comprises:

- Unimodal Encoders: Independently process image and text inputs, enabling standalone performance on vision-only or language-only tasks.
- Image-Grounded Text Encoder: Encodes text in the context of associated visual content, enhancing retrieval performance.

3.2 • Image-Grounded Decoder: Enables generation of text conditioned on image inputs, facilitating tasks like captioning, VQA, and dialog generation. This flexible architecture allows VLX to be adapted to a wide range of tasks without requiring modifications. The structural interaction encoder between the and decoder components is finely tuned to balance taskspecific while nuances maintaining generalization capabilities. This design also supports multitask learning, allowing the model to be simultaneously trained on diverse objectives, which further enhances its robustness and applicability.

3.3 Data Optimization Pipeline To combat noisy training data, we introduce a Captioning and Filtering Loop (CapFilt):

To address the issue of noisy training data, we have developed a data optimization pipeline called capfilt. This pipeline includes a captioning and filtering loop that helps improve the quality of the data.

- caption generator: trained on a clean subset, this component generates diverse synthetic captions for the entire dataset, improving semantic alignment and diversity
- noise filter: uses both rule-based heuristics and model confidence to filter out lowquality original and synthetic captions, enhancing the reliability of supervision signals
- self-improving iteration: the filtered dataset is used to retrain the generator, continuously improving caption quality in a feedback loop, mimicking a form of selfsupervised refinement

This loop greatly improves the effectiveness and stability of training without the need for manual data labeling. Moreover, it allows the creation of taskspecific fake labels, expanding the model's usefulness to new datasets and areas of study. The pipeline also enables domain adaptation, as it enables vlx to adapt to unfamiliar distributions by iteratively improving data quality based on the context and target domain.

IV Experiments and Results

4.1 Datasets We evaluate VLX across multiple standard benchmarks:

- MS-COCO: Used for image captioning and retrieval.
- Flickr30K: Evaluated for image-text retrieval.
- VQA v2.0: Assesses visual question answering capabilities.
- MSR-VTT: Employed for video-language zeroshot transfer evaluation.

Additional ablations are conducted on smaller datasets such as Visual Dialog and NoCaps to test generalization capabilities. These datasets span a range of tasks, domains,



and modalities, providing a rigorous evaluation of VLX's adaptability.

4.2 Implementation Details Vlx is implemented in PyTorch, with backbone encoders initialized from pretrained resnet and bert variants. The model is trained using a mix of contrastive loss, image-text matching loss, and language modeling objectives. To expedite convergence, mixed- precision training and large-batch optimizations are employed. The capfilt module is set up using captions created by blip and adjusted using vlx's own outputs. Experiments are conducted on high-performance GPUs to efficiently handle largescale distributed training across multiple nodes.

4.3 Quantitative Results

Task	Metric	Gain Over
		Baseline
Image-	Recall@1	+3.1%
Text		
Retrieval		
Visual	Accuracy	+2.0%
Question	-	
Answering		
Image	CIDEr	+2.5%
Captioning		
Video	Recall@1	+1.8%
Retrieval		
(zero-shot)		

These improvements validate both the architectural flexibility and the effectiveness of the data optimization strategy. In addition to raw accuracy gains, VLX demonstrates faster convergence rates and better generalization to unseen tasks.

4.4 Ablation Studies To evaluate the impact of various components:

- With vs. Without CapFilt: Performance drops notably without the filtering loop, highlighting its central role.
- Single vs. Diverse Captions: Models trained on multiple diverse captions outperform those with only one synthetic caption per image.
- End-to-End vs. Modular Training: Joint optimization of encoder and decoder yields better results than separate training.
- Cross-modal Consistency: Aligning visual and textual attention improves interpretability and performance.

These studies confirm that each design choice contributes meaningfully to overall performance, and that the datacentric enhancements provide measurable gains across both seen and unseen benchmarks.

V Discussion

Vlx tackles persistent challenges in the vision-language pre-training domain by combining a flexible model architecture with high-quality data. The main advantage of vlx is its ability to work effectively in various tasks, thanks to its flexible architecture that can adapt to different network structures without major modifications. This broad applicability is especially beneficial in real-world scenarios where systems must be flexible enough to handle a wide range of tasks and environments. One significant contribution is the implementation of the capfilt pipeline, which streamlines the previously time-consuming process of data cleaning. Through an iterative process of generating and refining data, vlx improves the semantic relevance and variety of training material, eliminating the need for manual annotations. This process allows the model to constantly update its training data, resulting in enhanced resilience and decreased vulnerability to overfitting. Additionally, vlx exhibits exceptional zero-shot capabilities, suggesting its ability to apply its knowledge to new and unfamiliar domains. This is a vital feature for deployment in situations where labeled data is limited or inaccessible, such as languages with low resources or specific industry applications. The consistent enhancements in performance across various benchmark tests validate the effectiveness of our design choices. While vlx demonstrates impressive outcomes, there are still opportunities for further enhancement. For example, expanding the scope of captioning to include multiple modes of communication, incorporating feedback from downstream tasks, and investigating curriculum learning approaches could enhance the overall performance. In a similar vein, studying vlx's behavior in diverse language and domain environments could pave the way for groundbreaking advancements in global and cross-cultural artificial intelligence systems.

VI Conclusion

We present vlx, a unified framework that aims to connect the understanding and generation tasks, bridging the gap between them. By combining a hybrid architecture and a new captioning-filtering loop, vlx tackles two significant challenges in vlp research: the lack of flexibility in models and the presence of data noise. Our experiments showcase that vlx achieves state-of-the-art performance on various benchmarks, indicating not only superior accuracy but also enhanced generalization in scenarios where no training data is available. The model's design prioritizes effective multitasking, and the self-improving data pipeline greatly reduces the reliance on human-labeled data. In addition to



these performance improvements, vlx introduces a flexible method for handling noisy data obtained from the web, which is often a significant obstacle in largescale vlp systems. The automated captioning and filtering loop (capfilt) guarantees that the training data remains coherent and contextually appropriate, leading to more effective and efficient learning. The iterative nature of capfilt not only enhances training stability but also establishes a foundation for ongoing learning and adaptation, which are essential for successful deployment ever-changing environments. Moreover, in vlx demonstrates strong transferability across modalities and tasks. Its ability to perform without any additional training on unseen data types and scenarios indicates its effectiveness in various situations. This makes it a suitable candidate for real-world applications where data availability, quality, and domain variability are common challenges. The framework's capability to effortlessly adjust across different visual, textual, and temporal modalities underscores its strength. Looking ahead, vlx offers a strong foundation for future research in several promising directions. These are some of the areas of research that are being explored, such as learning multiple languages, reasoning with different types of information in real-time, creating personalized content, and developing artificial intelligence that can understand and interact with humans in a social way. The modular structure of vlx makes it highly adaptable for integration with cutting-edge technologies like augmented reality, robotics, and conversational artificial intelligence systems. Ultimately, we believe vlx represents a significant step toward building general-purpose, intelligent models that can handle various types of multimodal data. Its unique combination of architectural adaptability, data-driven optimization, and performance versatility sets a new benchmark for the upcoming generation of vlp systems. We anticipate that the availability of our code, datasets, and pre-trained models will expedite advancements in multimodal AI and spark fresh ideas and breakthroughs in the field.

VII References

 Saklayen, M. G. The global epidemic of the metabolic syndrome. Current hypertension reports, 2018, 20(2), 1-8.
Moradi, M., & Ghadiri, N. Different approaches for identifying important concepts in probabilistic biomedical text summarization. Artificial intelligence in medicine, 2018, 84, 101-116.

[3] Moradi, M., & Samwald, M. Deep Learning, Natural Language Processing, and Explainable Artificial Intelligence in the Biomedical Domain. 2022, arXiv preprint arXiv:2202.12678. [4] Demsar J, Curk T, Erjavec A, Gorup C, Hocevar T, Milutinovic M, Mozina M, Polajnar M, Toplak M, Staric A, Stajdohar M, Umek L, Zagar L, Zbontar J, Zitnik M, Zupan B. Orange: Data Mining Toolbox in Python, Journal of Machine Learning Researchm, 2013. 14(Aug): 2349–2353.

[5] Xia, S., Zhang, J., Du, G., Li, S., Vong, C. T., Yang, Z., ... & Li, C. A microcosmic syndrome differentiation model for metabolic syndrome with multilabel learning. Evidence-Based Complementary and Alternative Medicine, 2020.

[6] Wang, H., Wang, Y., Li, X., Deng, X., Kong, Y., Wang, W., & Zhou, Y. Machine learning of plasma metabolome identifies biomarker panels for metabolic syndrome: findings from the China Suboptimal Health Cohort. Cardiovascular Diabetology. 2022.

[7] Y. Wang, D. Li, X. Xu, Q. Jia, Z. Yang, W. Nai, and Y. Sun, "Logistic regression with variable fractional gradient descent method," IEEE 9th International Information Technology and Artificial Intelligence Conference (ITAIC), 2020.

[8] Abro, Sindhu, et al. "Automatic hate speech detection using machine learning: A comparative study." Machine Learning 10.6 (2020).

[9] T. Inoue and S. Abe, "Fuzzy support vector machines for pattern classification", In Proceedings of International Joint Conference on Neural Networks (IJCNN'01), Vol. 2, pp 1449-1454, July 2001.

[10] K. Zhou, W. Nai, S. Zhu, S. Zhang, Y. Xing, Z. Yang, and D. Li, "Logistic regression based on bat-inspired algorithm with Gaussian initialization," IEEE 5th Advanced Information Technology, Electronic and Automation Control Conference (IAEAC), 2021.

[11] G. Heinze, and M. Schemper, "A solution to the problem of separation in logistic regression," Statistics in Medicine, vol. 21, no. 16, 2002.

[12] C. Mood, "Logistic regression: Why we cannot do what we think we can do, and what we can do about it," European Sociological Review, vol. 26, no. 1, 2010.

L